

QPLEX: DUPLEX DUELING MULTI-AGENT Q-LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

We explore value-based multi-agent reinforcement learning (MARL) in the popular paradigm of centralized training with decentralized execution (CTDE). CTDE has an important concept, *Individual-Global-Max* (IGM) principle, which requires the consistency between joint and local action selections to support efficient local decision-making. However, in order to achieve scalability, existing MARL methods either limit representation expressiveness of their value function classes or relax the IGM consistency, which may suffer from instability risk or lead to poor performance. This paper presents a novel MARL approach, called *duPLEX dueling multi-agent Q-learning* (QPLEX), which takes a duplex dueling network architecture to factorize the joint value function. This duplex dueling structure encodes the IGM principle into the neural network architecture and thus enables efficient value function learning. Theoretical analysis shows that QPLEX achieves a complete IGM function class. Empirical experiments on StarCraft II micromanagement tasks demonstrate that QPLEX significantly outperforms state-of-the-art baselines in both online and offline data collection settings, and also reveal that QPLEX achieves high sample efficiency and can benefit from offline datasets without additional online exploration.

1 INTRODUCTION

Cooperative multi-agent reinforcement learning (MARL) has broad prospects for addressing many complex real-world problems, such as sensor networks (Zhang & Lesser, 2011), coordination of robot swarms (Hüttenrauch et al., 2017), and autonomous cars (Cao et al., 2012). However, cooperative MARL encounters two major challenges of scalability and partial observability in practical applications. The joint state-action space grows exponentially as the number of agents increases. The partial observability and communication constraints of the environment require each agent to make its individual decisions based on local action-observation histories. To address these challenges, a popular MARL paradigm, called *centralized training with decentralized execution* (CTDE) (Oliehoek et al., 2008; Kraemer & Banerjee, 2016), has recently attracted great attention, where agents’ policies are trained with access to global information in a centralized way and executed only based on local histories in a decentralized way.

Many CTDE learning approaches have been proposed recently, among which value-based MARL algorithms (Sunehag et al., 2018; Rashid et al., 2018; Son et al., 2019; Yang et al., 2020) have shown state-of-the-art performance on challenging tasks, e.g., unit micromanagement in StarCraft II (Samvelyan et al., 2019). To enable effective CTDE for multi-agent Q-learning, it is critical that the joint greedy action should be equivalent to the collection of individual greedy actions of agents, which is called the IGM (*Individual-Global-Max*) principle (Son et al., 2019). This IGM principle provides two advantages: 1) ensuring the policy consistency during centralized training (learning the joint Q-function) and decentralized execution (using individual Q-functions) and 2) enabling scalable centralized training of computing one-step TD target of the joint Q-function (deriving joint greedy action selection from individual Q-functions). To enable this principle, VDN (Sunehag et al., 2018) and QMIX (Rashid et al., 2018) propose two sufficient conditions of IGM to factorize the joint action-value function. However, these two decomposition methods suffer from structural constraints and limit the joint action-value function class they can represent. As shown by Wang et al. (2020), the incompleteness of the joint value function class may lead to poor performance or potential risk of

training instability in the offline setting (Levine et al., 2020). To address this structural limitation, QTRAN (Son et al., 2019) proposes a factorization method expressing the complete value function space induced by the IGM consistency, but its exact implementation is known to be computationally intractable. The approximate version of QTRAN requires two extra soft regularizations and performs poorly in complex domains with online data collection (Mahajan et al., 2019). Therefore, achieving effective scalability remains an open problem for cooperative MARL.

To address this challenge, this paper presents a novel MARL approach, called *duPLEX dueling multi-agent Q-learning* (QPLEX), that takes a duplex dueling network architecture to factorize the joint action-value function into individual action-value functions. QPLEX introduces the dueling structure $Q = V + A$ (Wang et al., 2016) for representing both joint and individual (duplex) action-value functions and then reformalizes the IGM principle as an *advantage-based IGM*. This reformulation transforms the IGM consistency into the constraints on the value range of the advantage functions and thus facilitates the action-value function learning with linear decomposition structure. Unlike QTRAN that uses soft constraints and loses the guarantee of exact IGM consistency (Son et al., 2019), QPLEX takes advantage of a duplex dueling architecture to encode it into the neural network structure and provide a guaranteed IGM consistency. To our best knowledge, QPLEX is the first multi-agent Q-learning algorithm that effectively achieves high scalability with a full realization of the IGM principle.

We evaluate the performance of QPLEX in both didactic problems proposed by prior work (Son et al., 2019; Wang et al., 2020) and a range of unit micromanagement benchmark tasks in StarCraft II (Samvelyan et al., 2019). In these didactic problems, QPLEX demonstrates its full representation expressiveness, thereby learning the optimal policy and avoiding the potential risk of training instability. Empirical results on more challenging StarCraft II tasks show that QPLEX significantly outperforms other multi-agent Q-learning baselines in online and offline data collections. It is particularly interesting that QPLEX shows the ability to support offline training, which is not possessed by other baselines. This ability not only provides QPLEX with high stability and sample efficiency but also with opportunities to efficiently utilize multi-source offline data without additional online exploration (Fujimoto et al., 2019; Fu et al., 2020; Levine et al., 2020; Yu et al., 2020).

2 PRELIMINARIES

2.1 DECENTRALIZED PARTIALLY OBSERVABLE MDP (DEC-POMDP)

We model a fully cooperative multi-agent task as a Dec-POMDP (Oliehoek et al., 2016) defined by a tuple $\mathcal{M} = \langle \mathcal{N}, \mathcal{S}, \mathcal{A}, P, \Omega, O, r, \gamma \rangle$, where $\mathcal{N} \equiv \{1, 2, \dots, n\}$ is a finite set of agents and $s \in \mathcal{S}$ is a finite set of global states. At each time step, every agent $i \in \mathcal{N}$ chooses an action $a_i \in \mathcal{A} \equiv \{\mathcal{A}^{(1)}, \dots, \mathcal{A}^{(|\mathcal{A}|)}\}$ on a global state s , which forms a joint action $\mathbf{a} \equiv [a_i]_{i=1}^n \in \mathcal{A} \equiv \mathcal{A}^n$. It results in a joint reward $r(s, \mathbf{a})$ and a transition to the next global state $s' \sim P(\cdot | s, \mathbf{a})$. $\gamma \in [0, 1]$ is a discount factor. We consider a *partially observable* setting, where each agent i receives an individual partial observation $o_i \in \Omega$ according to the observation probability function $O(o_i | s, a_i)$. Each agent i has an action-observation history $\tau_i \in \mathcal{T} \equiv (\Omega \times \mathcal{A})^*$ and constructs its individual policy $\pi_i(a | \tau_i)$ to jointly maximize team performance. We use $\tau \in \mathcal{T} \equiv \mathcal{T}^n$ to denote joint action-observation history. The formal objective function is to find a joint policy $\pi = \langle \pi_1, \dots, \pi_n \rangle$ that maximizes a joint value function $V^\pi(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, \pi]$. Another quantity of interest in policy search is the joint action-value function $Q^\pi(s, \mathbf{a}) = r(s, \mathbf{a}) + \gamma \mathbb{E}_{s'}[V^\pi(s')]$.

2.2 DEEP MULTI-AGENT Q-LEARNING IN DEC-POMDP

Q-learning algorithms is a popular algorithm to find the optimal joint action-value function $Q^*(s, \mathbf{a}) = r(s, \mathbf{a}) + \gamma \mathbb{E}_{s'}[\max_{\mathbf{a}'} Q^*(s', \mathbf{a}')]$. Deep Q-learning represents the action-value function with a deep neural network parameterized by θ . Multi-agent Q-learning algorithms (Sunehag et al., 2018; Rashid et al., 2018; Son et al., 2019; Yang et al., 2020) use a replay memory D to store the transition tuple $(\tau, \mathbf{a}, r, \tau')$, where r is the reward for taking action \mathbf{a} at joint action-observation history τ with a transition to τ' . Due to partial observability, $Q(\tau, \mathbf{a}; \theta)$ is used in place of $Q(s, \mathbf{a}; \theta)$. Thus, parameters θ are learnt by minimizing the following expected TD error:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\tau, \mathbf{a}, r, \tau') \in D} \left[(r + \gamma V(\tau'; \theta^-) - Q(\tau, \mathbf{a}; \theta))^2 \right], \quad (1)$$

where $V(\tau'; \theta^-) = \max_{\mathbf{a}'} Q(\tau', \mathbf{a}'; \theta^-)$ is the one-step expected future return of the TD target and θ^- are the parameters of the target network, which will be periodically updated with θ .

2.3 CENTRALIZED TRAINING WITH DECENTRALIZED EXECUTION (CTDE)

CTDE is a popular paradigm of cooperative multi-agent deep reinforcement learning (Sunehag et al., 2018; Rashid et al., 2018; Son et al., 2019; Yang et al., 2020). Agents are trained in a centralized way and granted access to other agents' information or the global states during the centralized training process. However, due to partial observability and communication constraints, each agent makes its own decision based on its local action-observation history during the decentralized execution phase. IGM (*Individual-Global-Max*) (Son et al., 2019) is a popular principle to realize effective value-based CTDE, which asserts the consistency between joint and local greedy action selections in the joint action-value $Q_{tot}(\tau, \mathbf{a})$ and individual action-values $[Q_i(\tau_i, a_i)]_{i=1}^n$:

$$\forall \tau \in \mathcal{T}, \arg \max_{\mathbf{a} \in \mathcal{A}} Q_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} Q_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} Q_n(\tau_n, a_n) \right). \quad (2)$$

Two factorization structures, **additivity** and **monotonicity**, has been proposed by VDN (Sunehag et al., 2018) and QMIX (Rashid et al., 2018), respectively, as shown below:

$$Q_{tot}^{\text{VDN}}(\tau, \mathbf{a}) = \sum_{i=1}^n Q_i(\tau_i, a_i) \quad \text{and} \quad \forall i \in \mathcal{N}, \frac{\partial Q_{tot}^{\text{QMIX}}(\tau, \mathbf{a})}{\partial Q_i(\tau_i, a_i)} > 0. \quad (3)$$

These two structures are sufficient conditions for the IGM constraint. Qatten (Yang et al., 2020) is a variant of VDN, which supplements global information through a multi-head attention structure. However, they are not necessary conditions and limit the representation expressiveness of joint action-value functions (Mahajan et al., 2019). There exist tasks whose factorizable joint action-value functions can not be represented by these decomposition methods, as shown in Section 4. In contrast, QTRAN (Son et al., 2019) transforms IGM into a linear constraint and uses it as soft regularization constrains. However, this relaxation may violate the exact IGM consistency and result in poor performance in complex problems.

3 QPLEX: DUPLEX DUELING MULTI-AGENT Q-LEARNING

In this section, we will first introduce advantage-based IGM, equivalent to the regular IGM principle, and, with this new definition, convert the IGM consistency of greedy action selection to simple constraints on advantage functions. We then present a novel deep MARL model, called *duPLEX dueling multi-agent Q-learning algorithm* (QPLEX), that directly realizes these constraints by a scalable neural network architecture.

3.1 ADVANTAGE-BASED IGM

To ensure the consistency of greedy action selection on the joint and local action-value functions, the IGM principle constrains the relative order of Q-values over actions. From the perspective of dueling decomposition structure $Q = V + A$ proposed by Dueling DQN (Wang et al., 2016), this consistency should only constrain the action-dependent advantage term A and be free of the state-value function V . This observation naturally motivates us to reformalize the IGM principle as advantage-based IGM, which transforms the consistency constraint onto advantage functions.

Definition 1 (Advantage-based IGM). *For a joint action-value function $Q_{tot}: \mathcal{T} \times \mathcal{A} \mapsto \mathbb{R}$ and individual action-value functions $[Q_i: \mathcal{T} \times \mathcal{A} \mapsto \mathbb{R}]_{i=1}^n$, where $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{N}$,*

$$\textbf{(Joint Dueling)} \quad Q_{tot}(\tau, \mathbf{a}) = V_{tot}(\tau) + A_{tot}(\tau, \mathbf{a}) \text{ and } V_{tot}(\tau) = \max_{\mathbf{a}'} Q_{tot}(\tau, \mathbf{a}'), \quad (4)$$

$$\textbf{(Individual Dueling)} \quad Q_i(\tau_i, a_i) = V_i(\tau_i) + A_i(\tau_i, a_i) \text{ and } V_i(\tau_i) = \max_{a'_i} Q_i(\tau_i, a'_i), \quad (5)$$

such that the following holds

$$\arg \max_{\mathbf{a} \in \mathcal{A}} A_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} A_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} A_n(\tau_n, a_n) \right), \quad (6)$$

then, we can say that $[Q_i]_{i=1}^n$ satisfies advantage-based IGM for Q_{tot} .

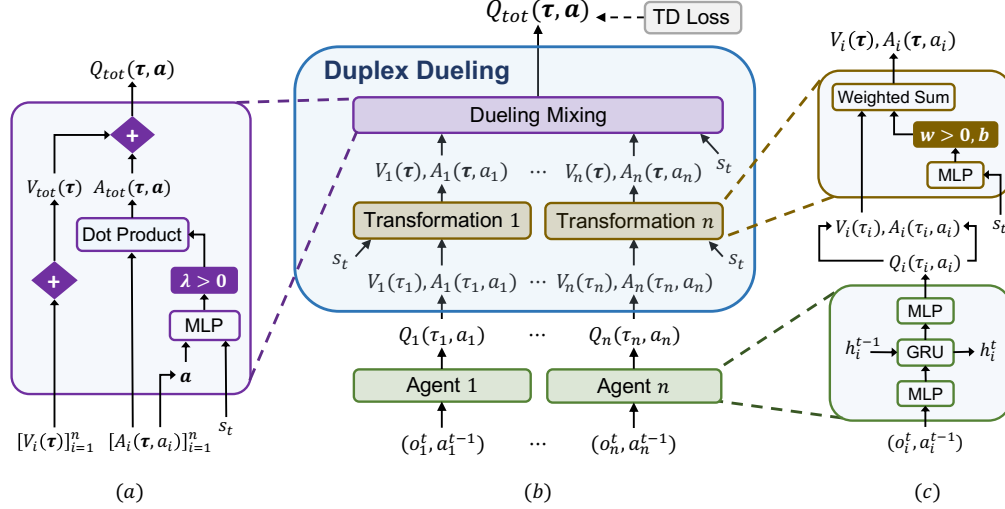


Figure 1: (a) The dueling mixing network structure. (b) The overall QPLEX architecture. (c) Agent network structure (bottom) and Transformation network structure (top).

As specified in Definition 1, advantage-based IGM takes a duplex dueling architecture, *Joint Dueling* and *Individual Dueling*, which induces the joint and local (duplex) advantage functions by $A = Q - V$. Compared with regular IGM, advantage-based IGM transfers the consistency constraint on action-value functions stated in Eq. (2) to that on advantage functions. This change is an equivalent transformation because the state-value terms V do not affect the action selection, as shown by Proposition 1.

Proposition 1. *The advantage-based IGM and IGM function classes are equivalent.*

One key benefit of using advantage-based IGM is that its consistency constraint can be directly realized by limiting the value range of advantage functions, as indicated by the following fact.

Fact 1. *The constraint of advantage-based IGM stated in Eq. (6) is equivalent to that when $\forall \tau \in \mathcal{T}$, $\forall \mathbf{a}^* \in \mathcal{A}^*(\tau)$, $\forall \mathbf{a} \in \mathcal{A} \setminus \mathcal{A}^*(\tau)$, $\forall i \in \mathcal{N}$,*

$$A_{tot}(\tau, \mathbf{a}^*) = A_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad A_{tot}(\tau, \mathbf{a}) < 0, A_i(\tau_i, a_i) \leq 0, \quad (7)$$

where $\mathcal{A}^*(\tau) = \{\mathbf{a} | \mathbf{a} \in \mathcal{A}, Q_{tot}(\tau, \mathbf{a}) = V_{tot}(\tau)\}$.

To achieve a full expressiveness power of advantage-based IGM or IGM, Fact 1 enables us to develop an efficient MARL algorithm that allows the joint state-value function learning with any scalable decomposition structure and just imposes simple constraints limiting value ranges of advantage functions. The next subsection will describe such a MARL algorithm.

3.2 THE QPLEX ARCHITECTURE

In this subsection, we present a novel multi-agent Q-learning algorithm with a duplex dueling architecture, called QPLEX, which exploits Fact 1 and realizes the advantage-based IGM constraint. The overall architecture of QPLEX is illustrated in Figure 1, which consists of two main components as follows: (i) an *Individual Action-Value Function* for each agent, and (ii) a *Duplex Dueling* component that composes individual action-value functions into a joint action-value function under the advantage-based IGM constraint. During the centralized training, the whole network is learned in an end-to-end fashion to minimize the TD loss as specified in Eq. (1). During the decentralized execution, the duplex dueling component will be removed and each agent will select actions using its individual Q-function based on local action-observation histories.

Individual Action-Value Function is represented by a recurrent Q-network for each agent i , which takes previous hidden state h_i^{t-1} , current local observations o_i^t , and previous action a_i^{t-1} as inputs and outputs local $Q_i(\tau_i, a_i)$.

Duplex Dueling component connects local and joint action-value functions via two modules: (i) a *Transformation* network module that incorporates the information of global state or joint history into individual action-value functions during the centralized training process, and (ii) a *Dueling Mixing* network module that composes separate action-value functions from *Transformation* into

a joint action-value function. *Duplex Dueling* first derives the individual dueling structure for each agent i by computing its value function $V_i(\tau_i) = \max_{a_i} Q_i(\tau_i, a_i)$ and its advantage function $A_i(\tau_i, a_i) = Q_i(\tau_i, a_i) - V_i(\tau_i)$, and then computes the joint dueling structure by using individual dueling structures.

Transformation network module uses the centralized information to transform local dueling structure $[V_i(\tau_i), A_i(\tau_i, a_i)]_{i=1}^n$ to $[V_i(\boldsymbol{\tau}), A_i(\boldsymbol{\tau}, a_i)]_{i=1}^n$ conditioned on the joint action-observation history, as shown below, for any agent i ,

$$V_i(\boldsymbol{\tau}) = w_i(\boldsymbol{\tau})V_i(\tau_i) + b_i(\boldsymbol{\tau}) \quad \text{and} \quad A_i(\boldsymbol{\tau}, a_i) = w_i(\boldsymbol{\tau})A_i(\tau_i, a_i) + b_i(\boldsymbol{\tau}), \quad (8)$$

where $w_i(\boldsymbol{\tau}) > 0$ is a positive weight. This positive linear transformation maintains the consistency of the greedy action selection and alleviates partial observability in Dec-POMDP (Son et al., 2019; Yang et al., 2020). As used by QMIX (Rashid et al., 2018), QTRAN (Son et al., 2019), and Qatten (Yang et al., 2020), the centralized information can be the global state s , if available, or the joint action-observation history $\boldsymbol{\tau}$.

Dueling Mixing network module takes the outputs of the transformation network as input, e.g., $[V_i, A_i]_{i=1}^n$, and produces the values of joint Q_{tot} , as shown in Figure 1a. This dueling mixing network uses individual dueling structure transformed by *Transformation* to compute the joint value $V_{tot}(\boldsymbol{\tau})$ and the joint advantage $A_{tot}(\boldsymbol{\tau}, \mathbf{a})$, respectively, and finally outputs $Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = V_{tot}(\boldsymbol{\tau}) + A_{tot}(\boldsymbol{\tau}, \mathbf{a})$ by using the joint dueling structure.

Based on Fact 1, the advantage-based IGM principle imposes no constraints on value functions. Therefore, to enable efficient learning, we use a simple sum structure to compose the joint value:

$$V_{tot}(\boldsymbol{\tau}) = \sum_{i=1}^n V_i(\boldsymbol{\tau}) \quad (9)$$

To enforce the IGM consistency of the joint advantage and individual advantages, as specified by Eq. (7), QPLEX computes the joint advantage function as follows:

$$A_{tot}(\boldsymbol{\tau}, \mathbf{a}) = \sum_{i=1}^n \lambda_i(\boldsymbol{\tau}, \mathbf{a}) A_i(\boldsymbol{\tau}, a_i), \quad \text{where } \lambda_i(\boldsymbol{\tau}, \mathbf{a}) > 0. \quad (10)$$

The joint advantage function A_{tot} is the dot product of separate advantage functions $[A_i]_{i=1}^n$ and positive importance weights $[\lambda_i]_{i=1}^n$. This positivity induced by λ_i will continue to maintain the consistency flow of the greedy action selection. To enable efficient learning of importance weights λ_i with joint history and action, QPLEX uses a scalable multi-head attention module (Vaswani et al., 2017):

$$\lambda_i(\boldsymbol{\tau}, \mathbf{a}) = \sum_{k=1}^K \lambda_{i,k}(\boldsymbol{\tau}, \mathbf{a}) \phi_{i,k}(\boldsymbol{\tau}) v_k(\boldsymbol{\tau}), \quad (11)$$

where K is the number of attention heads, $\lambda_{i,k}(\boldsymbol{\tau}, \mathbf{a})$ and $\phi_{i,k}(\boldsymbol{\tau})$ are attention weights activated by a sigmoid regularizer, and $v_k(\boldsymbol{\tau}) > 0$ is a positive key of each head. This sigmoid activation of λ_i brings sparsity to the credit assignment of the joint advantage function to individuals, which enables efficient multi-agent learning (Wang et al., 2019).

With Eq. (9) and (10), the joint action-value function Q_{tot} can be reformulated as follows:

$$Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = V_{tot}(\boldsymbol{\tau}) + A_{tot}(\boldsymbol{\tau}, \mathbf{a}) = \sum_{i=1}^n Q_i(\boldsymbol{\tau}, a_i) + \sum_{i=1}^n (\lambda_i(\boldsymbol{\tau}, \mathbf{a}) - 1) A_i(\boldsymbol{\tau}, a_i). \quad (12)$$

It can be seen that Q_{tot} consists of two terms. The first term is the sum of separate action-value functions $[Q_i]_{i=1}^n$, which is basically the joint action-value function Q_{tot}^{Qatten} of Qatten (Yang et al., 2020) (which is the Q_{tot} of VDN (Sunehag et al., 2018) with global information). The second term corrects for the discrepancy between the centralized joint action-value function and Q_{tot}^{Qatten} , which is the main contribution of QPLEX to realize the full expressiveness power of value factorization.

Proposition 2. *Given the universal function approximation of neural networks, the action-value function class that QPLEX can realize is equivalent to what is induced by the IGM principle.*

In practice, QPLEX can utilize common neural network structures (e.g., multi-head attention modules) to achieve superior performance by approximating the universal approximation theorem (Csaji et al., 2001). We will discuss the effects of QPLEX’s duplex dueling network with different configurations in Section 4.1. As introduced by Son et al. (2019); Wang et al. (2020), the completeness of value factorization is very critical for multi-agent Q-learning and we will illustrate the stability and state-of-the-art performance of QPLEX in online and offline data collections in the next section.

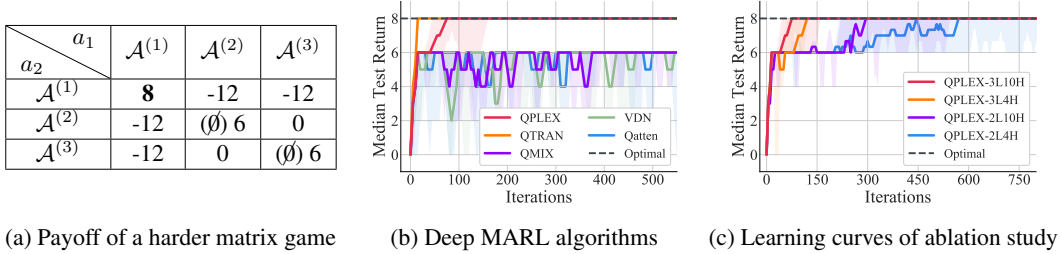


Figure 2: (a) Payoff matrix for a harder one-step game. Boldface means the optimal joint action selection from payoff matrix. (b) The learning curves of QPLEX and other baselines. (c) The learning curve of QPLEX, whose suffix $aLbH$ denotes the neural network size with a layers and b heads (multi-head attention) for learning importance weights λ_i (see Eq. (10) and (11)), respectively.

4 EXPERIMENTS

In this section, we first study didactic examples proposed by prior work (Son et al., 2019; Wang et al., 2020) to investigate the effects of QPLEX’s complete IGM expressiveness on learning optimality and stability. To demonstrate scalability on complex MARL domains, we also evaluate the performance of QPLEX on a range of StarCraft II benchmark tasks (Samvelyan et al., 2019). The completeness of the IGM function class can express richer joint action-value function classes induced by large and diverse datasets or training buffers. This expressiveness can provide QPLEX with higher sample efficiency to achieve state-of-the-art performance in online and offline data collections. We compare QPLEX with state-of-the-art baselines: QTRAN (Son et al., 2019), QMIX (Rashid et al., 2018), VDN (Sunehag et al., 2018), and Qatten (Yang et al., 2020). In particular, the second term of Eq. (12) is the main difference between QPLEX and Qatten. Thus, Qatten provides a natural ablation baseline of QPLEX to demonstrate the effectiveness of this discrepancy term. The implementation details of these algorithms and experimental settings are deferred to Appendix B. Towards fair evaluation, all experimental results are illustrated with the median performance and 25-75% percentiles over 6 random seeds. The videos of our experiments on StarCraft II are available on an anonymous website¹.

4.1 MATRIX GAMES

QTRAN (Son et al., 2019) proposes a hard matrix game, as shown in Table 4a of Appendix C. In this subsection, we consider a harder matrix game in Table 2a, which also describes a simple cooperative multi-agent task with considerable miscoordination penalties, and its local optimum is more difficult to jump out. The optimal joint strategy of these two games is to perform action $\mathcal{A}^{(1)}$ simultaneously. To ensure sufficient data collection in the joint action space, we adopt uniform data distribution. With this fixed dataset, we can study the optimality of multi-agent Q-learning from an optimization perspective, ignoring the challenge of exploration and sample complexity.

As shown in Figure 2b, QPLEX and QTRAN, which possess a richer expressiveness power of value factorization can achieve optimal performance, while other algorithms with limited expressiveness (e.g., QMIX, VDN, and Qatten) fall into a local optimum induced by miscoordination penalties. In the original matrix proposed by QTRAN, QPLEX and QTRAN are still the only two algorithms that can successfully converge to optimal joint action-value functions. These results are deferred to Appendix C. QTRAN achieves superior performance in the matrix games but suffers from its relaxation of IGM consistency in complex domains (such as StarCraft II) shown in Section 4.3.

In the theoretical analysis of QPLEX, Proposition 2 exploits the universal function approximation of neural networks. QPLEX allows scalable implementations with various neural network capacities (different layers and heads of attention module) for learning importance weights λ_i (see Eq. (10) and (11)). As shown in Figure 2c, by increasing the neural network size for learning λ_i (e.g., QPLEX-3L10H), QPLEX possesses more expressiveness of value factorization and converges faster. However, learning efficiency becomes challenging for complex neural networks. To effectively perform StarCraft II tasks ranging from 2 to 15 agents, we use a small multi-head attention module (i.e., QPLEX-1L4H) in complex domains (see Section 4.3). Please refer to Appendix B for more detailed configurations.

¹<https://sites.google.com/view/qplex-marl/>

4.2 TWO-STATE MMDP

In this subsection, we focus on a Multi-agent Markov Decision Process (MMDP) (Boutilier, 1996) which is a fully cooperative multi-agent setting with full observability. Consider a two-state MMDP proposed by Wang et al. (2020) with two agents, two actions, and a single reward (see Figure 6b of Appendix C). Two agents start at state s_2 and explore extrinsic rewards for 100 environment steps. The optimal policy of this MMDP is simply executing the action $\mathcal{A}^{(1)}$ at state s_2 , which is the only coordination pattern to obtain the positive reward. To approximate the uniform data distribution, we adopt a uniform exploration strategy (*i.e.*, ϵ -greedy exploration with $\epsilon = 1$). We consider the training stability of multi-agent Q-learning algorithms with uniform data distribution in this special MMDP task. As shown in Figure 3, the joint state-value function V_{tot} learned by baseline algorithms using limited function classes, including QMIX, VDN, and Qatten, will diverge. This instability phenomenon of VDN has been theoretically investigated by Wang et al. (2020). By utilizing richer function classes, both QPLEX and QTRAN can address this numerical instability issue and converge to the optimal joint state-value function.

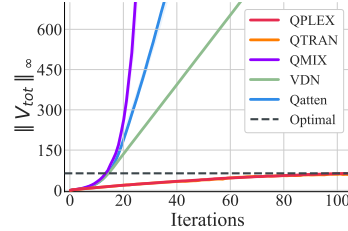


Figure 3: The learning curves of $\|V_{tot}\|_\infty$ in a specific two-state MMDP, which is shown in Figure 6b of Appendix C.

4.3 DECENTRALIZED STARCRAFT II MICROMANAGEMENT BENCHMARK

A more challenging set of empirical experiments are based on StarCraft Multi-Agent Challenge (SMAC) benchmark (Samvelyan et al., 2019). We first investigate empirical performance in a popular experimental setting with ϵ -greedy exploration and a limited first-in-first-out (FIFO) buffer (Samvelyan et al., 2019), named online data collection setting. To demonstrate the offline training potential of QPLEX, we also adopt the offline data collection setting proposed by Levine et al. (2020), which can be granted access to a given dataset without additional online exploration.

4.3.1 TRAINING WITH ONLINE DATA COLLECTION

Figure 4 shows the results of StarCraft II under the online data collection process, in which QPLEX significantly outperforms other baselines with higher sample efficiency. On the super hard map 5s10z, the performance gap between QPLEX and other baselines exceeds 30% in win rate, and the visualized strategies of QPLEX and QMIX in this map are deferred to Appendix B. Most multi-agent Q-learning baselines including QMIX, VDN, and Qatten achieve reasonable performance (see Figure 4). However, QTRAN performs the worst in these comparative experiments, even though it performs well in the didactic games. From a theoretical perspective, the online data collection process utilizes an ϵ -greedy exploration process, which requires individual greedy action selections to build an effective training buffer. QTRAN may suffer from its relaxation of IGM consistency (soft constraints of IGM) in the online data collection phase, while the duplex dueling architecture of QPLEX (hard constraint of IGM) provides effective individual greedy action selections, making it suitable for data collection with ϵ -greedy exploration.

4.3.2 TRAINING WITH OFFLINE DATA COLLECTION

Recently, offline reinforcement learning has been regarded as a key step for real-world RL applications (Dulac-Arnold et al., 2019; Levine et al., 2020). Agarwal et al. (2020) presents an optimistic perspective of offline Q-learning that DQN and its variants can achieve superior performance in Atari 2600 games (Bellemare et al., 2013) with sufficiently large and diverse datasets. In MARL, StarCraft II benchmark has the same discrete action space as Atari. We conduct a lot of experiments on the StarCraft II benchmark tasks to study offline multi-agent Q-learning in this subsection. Different from other related works that study *distributional shift* (Fujimoto et al., 2019; Levine et al., 2020), we adopt a large and diverse dataset to make the expressiveness power of value factorization become the dominant factor to investigate. We train a behavior policy of QMIX and collect all its experienced transitions throughout the training process (see the details in Appendix C). As shown in Figure 5, QPLEX significantly outperforms other multi-agent Q-learning baselines and possesses the state-of-the-art value factorization structure for offline multi-agent Q-learning. QMIX and Qatten cannot

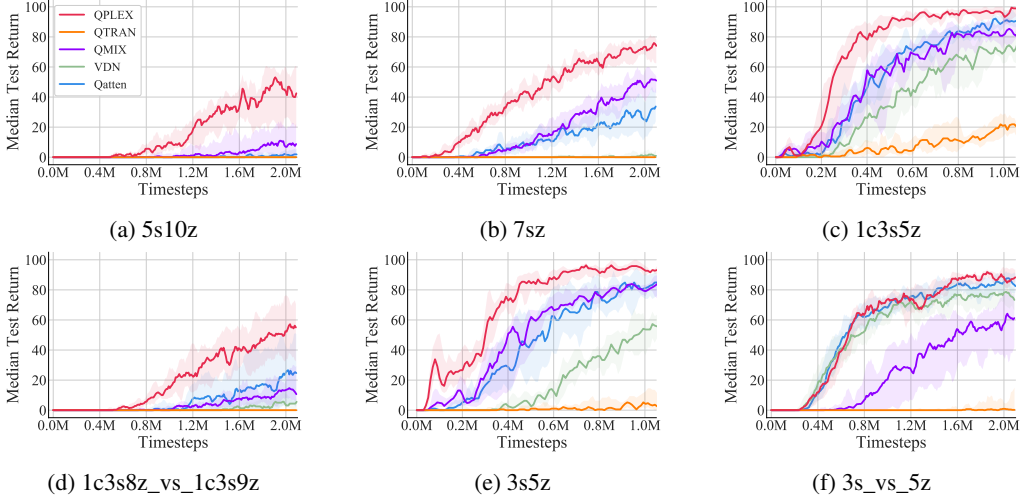


Figure 4: Learning curves of StarCraft II with online data collection on six different maps.

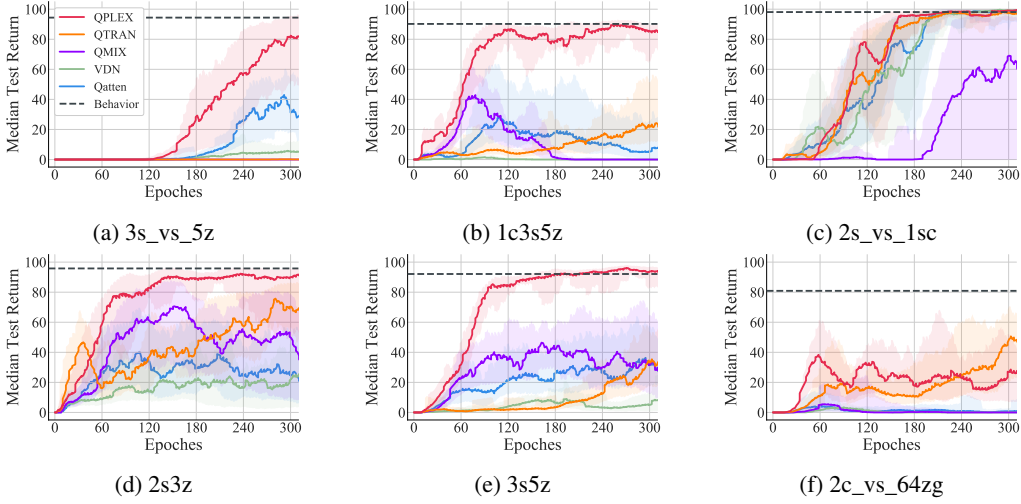


Figure 5: Learning curves of StarCraft II with offline data collection on six different maps.

always maintain stable learning performance, and VDN suffers from offline data collection and leads to weak empirical results. QTRAN may perform well in certain cases when its soft constraints, two ℓ_2 -penalty terms, are well minimized. With offline data collection, individual greedy action selections do not need to build a training buffer, but they still need to compute the one-step TD target for centralized training. Therefore, compared with QTRAN, QPLEX still has theoretical advantages regarding the IGM principle in the offline data collection setting.

5 CONCLUSION

In this paper, we introduced QPLEX, a novel multi-agent Q-learning framework that allows centralized end-to-end training and learns to factorize a joint action-value function to enable decentralized execution. QPLEX takes advantage of a duplex dueling architecture that efficiently encodes the IGM consistency constraint on joint and individual greedy action selections. Our theoretical analysis shows that QPLEX achieves a complete IGM function class. Empirical results demonstrate that it significantly outperforms state-of-the-art baselines in both online and offline data collection settings. In particular, QPLEX possesses a strong ability of supporting offline training. This ability provides QPLEX with high sample efficiency and opportunities of utilizing offline multi-source datasets. It will be an interesting and valuable direction to study offline multi-agent reinforcement learning in continuous action spaces (such as MuJoCo (Todorov et al., 2012)) with QPLEX’s value factorization.

REFERENCES

- Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, 2020.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47: 253–279, 2013.
- Craig Boutilier. Planning, learning and coordination in multiagent decision processes. In *Proceedings of the 6th Conference on Theoretical Aspects of Rationality and Knowledge*, pp. 195–210. Morgan Kaufmann Publishers Inc., 1996.
- Yongcan Cao, Wenwu Yu, Wei Ren, and Guanrong Chen. An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial Informatics*, 9(1): 427–438, 2012.
- Balázs Csanád Csáji et al. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Loránd University, Hungary*, 24(48):7, 2001.
- Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforcement learning. *arXiv preprint arXiv:1904.12901*, 2019.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pp. 2052–2062, 2019.
- Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. Guided deep reinforcement learning for swarm systems. *arXiv preprint arXiv:1709.06011*, 2017.
- Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, pp. 7611–7622, 2019.
- Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.
- Frans A Oliehoek, Christopher Amato, et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 4292–4301, 2018.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2186–2188. International Foundation for Autonomous Agents and Multiagent Systems, 2019.
- Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 5887–5896, 2019.

- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2085–2087, 2018.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Jianhao Wang, Zhizhou Ren, Beining Han, and Chongjie Zhang. Towards understanding linear value decomposition in cooperative multi-agent q-learning. *arXiv preprint arXiv:2006.00587*, 2020.
- Tonghan Wang, Jianhao Wang, Chongyi Zheng, and Chongjie Zhang. Learning nearly decomposable value functions via communication minimization. *arXiv preprint arXiv:1910.05366*, 2019.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International Conference on Machine Learning*, pp. 1995–2003, 2016.
- Yaodong Yang, Jianye Hao, Ben Liao, Kun Shao, Guangyong Chen, Wulong Liu, and Hongyao Tang. Qatten: A general framework for cooperative multiagent reinforcement learning. *arXiv preprint arXiv:2002.03939*, 2020.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *arXiv preprint arXiv:2005.13239*, 2020.
- Chongjie Zhang and Victor Lesser. Coordinated multi-agent reinforcement learning in networked distributed pomdps. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.

A OMITTED PROOFS IN SECTION 3

Definition 1 (Advantage-based IGM). *For a joint action-value function $Q_{tot}: \mathcal{T} \times \mathcal{A} \mapsto \mathbb{R}$ and individual action-value functions $[Q_i: \mathcal{T} \times \mathcal{A} \mapsto \mathbb{R}]_{i=1}^n$, where $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{N}$,*

$$\textbf{(Joint Dueling)} \quad Q_{tot}(\tau, \mathbf{a}) = V_{tot}(\tau) + A_{tot}(\tau, \mathbf{a}) \text{ and } V_{tot}(\tau) = \max_{\mathbf{a}'} Q_{tot}(\tau, \mathbf{a}'), \quad (4)$$

$$\textbf{(Individual Dueling)} \quad Q_i(\tau_i, a_i) = V_i(\tau_i) + A_i(\tau_i, a_i) \text{ and } V_i(\tau_i) = \max_{a_i'} Q_i(\tau_i, a_i'), \quad (5)$$

such that the following holds

$$\arg \max_{\mathbf{a} \in \mathcal{A}} A_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} A_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} A_n(\tau_n, a_n) \right), \quad (6)$$

then, we can say that $[Q_i]_{i=1}^n$ satisfies advantage-based IGM for Q_{tot} .

Let the action-value function class derived from IGM is denoted by

$$\tilde{\mathcal{Q}} = \left\{ \left(\tilde{Q}_{tot} \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|^n}, [\tilde{Q}_i \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|}]_{i=1}^n \right) \mid \text{Eq. (2) is satisfied} \right\}, \quad (13)$$

where \tilde{Q}_{tot} and $[\tilde{Q}_i]_{i=1}^n$ denote the joint and individual action-value functions induced by IGM, respectively. Similarly, let

$$\hat{\mathcal{Q}} = \left\{ \left(\hat{Q}_{tot} \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|^n}, [\hat{Q}_i \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|}]_{i=1}^n \right) \mid \text{Eq. (4), (5), (6) are satisfied} \right\} \quad (14)$$

denote the action-value function class derived from advantage-based IGM. \tilde{V}_{tot} and \tilde{A}_{tot} denote the joint state-value and advantage functions, respectively. $[\tilde{V}_i]_{i=1}^n$ and $[\tilde{A}_i]_{i=1}^n$ denote the individual state-value and advantage functions induced by advantage-IGM, respectively. According to the duplex dueling architecture $Q = V + A$ stated in advantage-based IGM (see Definition 1), we derive the joint and individual action-value functions as following: $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{N}$,

$$\hat{Q}_{tot}(\tau, \mathbf{a}) = \hat{V}_{tot}(\tau) + \hat{A}_{tot}(\tau, \mathbf{a}) \quad \text{and} \quad \hat{Q}_i(\tau_i, a_i) = \hat{V}_i(\tau_i) + \hat{A}_i(\tau_i, a_i). \quad (15)$$

Proposition 1. *The advantage-based IGM and IGM function classes are equivalent.*

Proof. We will prove $\tilde{\mathcal{Q}} \equiv \hat{\mathcal{Q}}$ in the following two directions.

$\tilde{\mathcal{Q}} \subseteq \hat{\mathcal{Q}}$ For any $(\tilde{Q}_{tot}, [\tilde{Q}_i]_{i=1}^n) \in \tilde{\mathcal{Q}}$, we construct $\hat{Q}_{tot} = \tilde{Q}_{tot}$ and $[\hat{Q}_i]_{i=1}^n = [\tilde{Q}_i]_{i=1}^n$. The joint and individual state-value/advantage functions induced by advantage-IGM

$$\hat{V}_{tot}(\tau) = \max_{\mathbf{a}'} \hat{Q}_{tot}(\tau, \mathbf{a}') \quad \text{and} \quad \hat{A}_{tot}(\tau, \mathbf{a}) = \hat{Q}_{tot}(\tau, \mathbf{a}) - \hat{V}_{tot}(\tau), \quad (16)$$

$$\hat{V}_i(\tau_i) = \max_{a_i'} \hat{Q}_i(\tau_i, a_i') \quad \text{and} \quad \hat{A}_i(\tau_i, a_i) = \hat{Q}_i(\tau_i, a_i) - \hat{V}_i(\tau_i), \quad \forall i \in \mathcal{N}, \quad (17)$$

are derived by Eq. (4) and Eq. (5), respectively. Because state-value functions do not affect the greedy action selection, $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}$,

$$\arg \max_{\mathbf{a} \in \mathcal{A}} \tilde{Q}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \tilde{Q}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \tilde{Q}_n(\tau_n, a_n) \right) \quad (18)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \hat{Q}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \hat{Q}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \hat{Q}_n(\tau_n, a_n) \right) \quad (19)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \left(\hat{Q}_{tot}(\tau, \mathbf{a}) - \hat{V}_{tot}(\tau) \right) = \quad (20)$$

$$\left(\arg \max_{a_1 \in \mathcal{A}} \left(\hat{Q}_1(\tau_1, a_1) - \hat{V}_1(\tau_1) \right), \dots, \arg \max_{a_n \in \mathcal{A}} \left(\hat{Q}_n(\tau_n, a_n) - \hat{V}_n(\tau_n) \right) \right) \quad (21)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \hat{A}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \hat{A}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \hat{A}_n(\tau_n, a_n) \right). \quad (22)$$

$$(23)$$

Thus, $(\hat{Q}_{tot}, [\hat{Q}_i]_{i=1}^n) \in \hat{\mathcal{Q}}$, which means that $\tilde{\mathcal{Q}} \subseteq \hat{\mathcal{Q}}$.

$\widehat{\mathcal{Q}} \subseteq \widetilde{\mathcal{Q}}$ We will prove this direction in the same way. For any $(\widehat{Q}_{tot}, [\widehat{Q}_i]_{i=1}^n) \in \widehat{\mathcal{Q}}$, we construct $\widetilde{Q}_{tot} = \widehat{Q}_{tot}$ and $[\widetilde{Q}_i]_{i=1}^n = [\widehat{Q}_i]_{i=1}^n$. Because state-value functions do not affect the greedy action selection, $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}$,

$$\arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{A}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \widehat{A}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \widehat{A}_n(\tau_n, a_n) \right) \quad (24)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \left(\widehat{A}_{tot}(\tau, \mathbf{a}) + \widehat{V}_{tot}(\tau) \right) = \quad (25)$$

$$\left(\arg \max_{a_1 \in \mathcal{A}} \left(\widehat{A}_1(\tau_1, a_1) + \widehat{V}_1(\tau_1) \right), \dots, \arg \max_{a_n \in \mathcal{A}} \left(\widehat{A}_n(\tau_n, a_n) + \widehat{V}_n(\tau_n) \right) \right) \quad (26)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{Q}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \widehat{Q}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \widehat{Q}_n(\tau_n, a_n) \right) \quad (27)$$

$$\Rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} \widetilde{Q}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \widetilde{Q}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \widetilde{Q}_n(\tau_n, a_n) \right). \quad (28)$$

$$(29)$$

Thus, $(\widetilde{Q}_{tot}, [\widetilde{Q}_i]_{i=1}^n) \in \widetilde{\mathcal{Q}}$, which means that $\widehat{\mathcal{Q}} \subseteq \widetilde{\mathcal{Q}}$. The action-value function classes derived from advantage-based IGM and IGM are equivalent. \square

Fact 1. The constraint of advantage-based IGM stated in Eq. (6) is equivalent to that when $\forall \tau \in \mathcal{T}$, $\forall \mathbf{a}^* \in \mathcal{A}^*(\tau)$, $\forall \mathbf{a} \in \mathcal{A} \setminus \mathcal{A}^*(\tau)$, $\forall i \in \mathcal{N}$,

$$A_{tot}(\tau, \mathbf{a}^*) = A_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad A_{tot}(\tau, \mathbf{a}) < 0, A_i(\tau_i, a_i) \leq 0, \quad (7)$$

where $\mathcal{A}^*(\tau) = \{\mathbf{a} | \mathbf{a} \in \mathcal{A}, Q_{tot}(\tau, \mathbf{a}) = V_{tot}(\tau)\}$.

Proof. We derive that $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{N}, \widehat{A}_{tot}(\tau, \mathbf{a}) \leq 0$ and $\widehat{A}_i(\tau_i, a_i) \leq 0$ from Eq. (4) and Eq. (5) of Definition 1, respectively. According to the definition of arg max operator, Eq. (4), and Eq. (5), $\forall \tau \in \mathcal{T}$, let $\widehat{\mathcal{A}}^*(\tau)$ denote $\arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{A}_{tot}(\tau, \mathbf{a})$ as follows:

$$\widehat{\mathcal{A}}^*(\tau) = \arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{A}_{tot}(\tau, \mathbf{a}) = \arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{Q}_{tot}(\tau, \mathbf{a}) \quad (30)$$

$$= \left\{ \mathbf{a} | \mathbf{a} \in \mathcal{A}, \widehat{Q}_{tot}(\tau, \mathbf{a}) = \widehat{V}_{tot}(\tau) \right\} \quad (31)$$

$$= \left\{ \mathbf{a} | \mathbf{a} \in \mathcal{A}, \widehat{Q}_{tot}(\tau, \mathbf{a}) - \widehat{V}_{tot}(\tau) = 0 \right\} \quad (32)$$

$$= \left\{ \mathbf{a} | \mathbf{a} \in \mathcal{A}, \widehat{A}_{tot}(\tau, \mathbf{a}) = 0 \right\}. \quad (33)$$

Similarly, $\forall \tau \in \mathcal{T}, \forall i \in \mathcal{N}$, let $\widehat{\mathcal{A}}_i^*(\tau_i)$ denote $\arg \max_{a_i \in \mathcal{A}} \widehat{A}_i(\tau_i, a_i)$ as follows:

$$\widehat{\mathcal{A}}_i^*(\tau_i) = \arg \max_{a_i \in \mathcal{A}} \widehat{A}_i(\tau_i, a_i) = \arg \max_{a_i \in \mathcal{A}} \widehat{Q}_i(\tau_i, a_i) \quad (34)$$

$$= \left\{ a_i | a_i \in \mathcal{A}, \widehat{Q}_i(\tau_i, a_i) = \widehat{V}_i(\tau_i) \right\} \quad (35)$$

$$= \left\{ a_i | a_i \in \mathcal{A}, \widehat{A}_i(\tau_i, a_i) = 0 \right\}. \quad (36)$$

Thus, $\forall \tau \in \mathcal{T}, \forall \mathbf{a}^* \in \widehat{\mathcal{A}}^*(\tau), \forall \mathbf{a} \in \mathcal{A} \setminus \widehat{\mathcal{A}}^*(\tau)$,

$$\widehat{A}_{tot}(\tau, \mathbf{a}^*) = 0 \quad \text{and} \quad \widehat{A}_{tot}(\tau, \mathbf{a}) < 0; \quad (37)$$

$\forall \tau \in \mathcal{T}, \forall i \in \mathcal{N}, \forall a_i^* \in \widehat{\mathcal{A}}_i^*(\tau_i), \forall a_i \in \mathcal{A} \setminus \widehat{\mathcal{A}}_i^*(\tau_i)$,

$$\widehat{A}_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad \widehat{A}_i(\tau_i, a_i) < 0. \quad (38)$$

Recall the constraint stated in Eq. 6, $\forall \tau \in \mathcal{T}$,

$$\arg \max_{\mathbf{a} \in \mathcal{A}} \widehat{A}_{tot}(\tau, \mathbf{a}) = \left(\arg \max_{a_1 \in \mathcal{A}} \widehat{A}_1(\tau_1, a_1), \dots, \arg \max_{a_n \in \mathcal{A}} \widehat{A}_n(\tau_n, a_n) \right). \quad (39)$$

We can rewrite the constraint of advantage-based IGM stated in Eq. (6) as $\forall \tau \in \mathcal{T}$,

$$\hat{\mathcal{A}}^*(\tau) = \left\{ \langle a_1, \dots, a_n \rangle \mid a_i \in \hat{\mathcal{A}}_i^*(\tau_i), \forall i \in \mathcal{N} \right\}. \quad (40)$$

Therefore, combining Eq. (37), Eq. (38), and Eq. (40), we can derive $\forall \tau \in \mathcal{T}, \forall \mathbf{a}^* \in \hat{\mathcal{A}}^*(\tau), \forall \mathbf{a} \in \mathcal{A} \setminus \hat{\mathcal{A}}^*(\tau), \forall i \in \mathcal{N}$,

$$\hat{A}_{tot}(\tau, \mathbf{a}^*) = \hat{A}_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad \hat{A}_{tot}(\tau, \mathbf{a}) < 0, \hat{A}_i(\tau_i, a_i) \leq 0. \quad (41)$$

In another way, combining Eq. (37), Eq. (38), and Eq. (41), we can derive Eq. (40) by the definition of $\hat{\mathcal{A}}^*$ and $[\hat{\mathcal{A}}^*]_{i=1}^n$ (see Eq. (33) and Eq. (36)). In more detail, the closed set property of Cartesian product of $[a_i^*]_{i=1}^n$ has been encoded into the Eq. (40) and Eq. (41) simultaneously. \square

Proposition 2. *Given the universal function approximation of neural networks, the action-value function class that QPLEX can realize is equivalent to what is induced by the IGM principle.*

Proof. We assume that the neural network of QPLEX can be large enough to achieve the universal function approximation by corresponding theorem (Csáji et al., 2001). Let the action-value function class that QPLEX can realize is denoted by

$$\overline{\mathcal{Q}} = \left\{ \left(\overline{Q}_{tot} \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|^n}, [\overline{Q}_i \in \mathbb{R}^{|\mathcal{T}||\mathcal{A}|}]_{i=1}^n \right) \mid \text{Eq. (8), (9), (10), (11), (12) are satisfied} \right\}. \quad (42)$$

In addition, $\overline{Q}_{tot}, \overline{V}_{tot}, \overline{A}_{tot}, [\overline{Q}'_i]_{i=1}^n, [\overline{V}'_i]_{i=1}^n, [\overline{A}'_i]_{i=1}^n, [\overline{Q}_i]_{i=1}^n, [\overline{V}_i]_{i=1}^n$, and $[\overline{A}_i]_{i=1}^n$ denote the corresponding (joint, transformed, and individual) (action-value, state-value, and advantage) functions, respectively. In the implementation of QPLEX, we ensure the positivity of important weights of *Transformation* and joint advantage function, $[w_i]_{i=1}^n$ and $[\lambda_i]_{i=1}^n$, which maintains the greedy action selection flow and rules out these non-interesting points (zeros) on optimization. We will prove $\hat{\mathcal{Q}} \equiv \overline{\mathcal{Q}}$ in the following two directions.

$\hat{\mathcal{Q}} \subseteq \overline{\mathcal{Q}}$ For any $(\hat{Q}_{tot}, [\hat{Q}_i]_{i=1}^n) \in \hat{\mathcal{Q}}$, we construct $\overline{Q}_{tot} = \hat{Q}_{tot}$ and $[\overline{Q}_i]_{i=1}^n = [\hat{Q}_i]_{i=1}^n$ and derive $\overline{V}_{tot}, \overline{A}_{tot}, [\overline{V}_i]_{i=1}^n$, and $[\overline{A}_i]_{i=1}^n$ by Eq.(4) and Eq. (5), respectively. In addition, we construct transformed functions connecting joint and individual functions as follows: $\forall \tau \in \mathcal{T}, \forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{N}$,

$$\overline{Q}'_i(\tau, \mathbf{a}) = \frac{\overline{Q}_{tot}(\tau, \mathbf{a})}{n}, \quad \overline{V}'_i(\tau) = \arg \max_{\mathbf{a} \in \mathcal{A}} \overline{Q}'_i(\tau, \mathbf{a}), \quad \text{and} \quad \overline{A}'_i(\tau, \mathbf{a}) = \overline{Q}'_i(\tau, \mathbf{a}) - \overline{V}'_i(\tau), \quad (43)$$

which means that according to Fact 1,

$$w_i(\tau) = 1, \quad b_i(\tau) = \overline{V}'_i(\tau) - \overline{V}_i(\tau_i), \quad \text{and} \quad \lambda_i(\tau, \mathbf{a}) = \begin{cases} \frac{\overline{A}'_i(\tau, \mathbf{a})}{\overline{A}_i(\tau_i, a_i)} > 0, & \text{when } \overline{A}_i(\tau_i, a_i) < 0, \\ 1, & \text{when } \overline{A}_i(\tau_i, a_i) = 0. \end{cases} \quad (44)$$

Thus, $(\overline{Q}_{tot}, [\overline{Q}_i]_{i=1}^n) \in \overline{\mathcal{Q}}$, which means that $\hat{\mathcal{Q}} \subseteq \overline{\mathcal{Q}}$.

$\overline{\mathcal{Q}} \subseteq \hat{\mathcal{Q}}$ For any $(\overline{Q}_{tot}, [\overline{Q}_i]_{i=1}^n) \in \overline{\mathcal{Q}}$, with the similar discussion of Fact 1, $\forall \tau \in \mathcal{T}, \forall i \in \mathcal{N}$, let $\overline{\mathcal{A}}_i^*(\tau_i)$ denote $\arg \max_{a_i \in \mathcal{A}} \overline{A}_i(\tau_i, a_i)$, where

$$\overline{\mathcal{A}}_i^*(\tau_i) = \{a_i \mid a_i \in \mathcal{A}, \overline{A}_i(\tau_i, a_i) = 0\}. \quad (45)$$

Combining the positivity of $[w_i]_{i=1}^n$ and $[\lambda_i]_{i=1}^n$ with Eq. (8), (9), (10), and (12), we can derive $\forall \tau \in \mathcal{T}, \forall i \in \mathcal{N}, \forall a_i^* \in \overline{\mathcal{A}}_i^*(\tau_i), \forall a_i \in \mathcal{A} \setminus \overline{\mathcal{A}}_i^*(\tau_i)$,

$$\overline{A}_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad \overline{A}_i(\tau_i, a_i) < 0 \quad (46)$$

$$\Rightarrow \overline{A}'_i(\tau, a_i^*) = w_i(\tau) \overline{A}_i(\tau_i, a_i^*) = 0 \quad \text{and} \quad \overline{A}'_i(\tau, a_i) = w_i(\tau) \overline{A}_i(\tau_i, a_i) < 0 \quad (47)$$

$$\Rightarrow \overline{A}_{tot}(\tau, \mathbf{a}^*) = \lambda_i(\tau, \mathbf{a}^*) \overline{A}'_i(\tau, a_i^*) = 0 \quad \text{and} \quad \overline{A}_{tot}(\tau, \mathbf{a}) = \lambda_i(\tau, \mathbf{a}) \overline{A}'_i(\tau, a_i) < 0, \quad (48)$$

QPLEX’s architecture configurations	Didactic Examples	StarCraft II
The number of layers in w, b, λ, ϕ, v	2 or 3	1
The number of heads in the attention module	4 or 10	4
Unit number in middle layers of w, b, λ, ϕ, v	64	\emptyset
Activation in the middle layers of w, v	Relu	\emptyset
Activation in the last layer of w, v	Absolute	Absolute
Activation in the middle layers of b	Relu	\emptyset
Activation in the last layer of b	None	None
Activation in the middle layers of λ, ϕ	Relu	\emptyset
Activation in the last layer of λ, ϕ	Sigmoid	Sigmoid

Table 1: The network configurations of QPLEX’s architecture.

where $\mathbf{a}^* = \langle a_1^*, \dots, a_n^* \rangle$ and $\mathbf{a} = \langle a_1, \dots, a_n \rangle$. Notably, these \mathbf{a}^* forms

$$\overline{\mathcal{A}}^*(\tau) = \left\{ \langle a_1, \dots, a_n \rangle \mid a_i \in \overline{\mathcal{A}}_i^*(\tau_i), \forall i \in \mathcal{N} \right\} \quad (49)$$

which is similar to Eq. (40) in the proof of Fact 1. We construct $\widehat{Q}_{tot} = \overline{Q}_{tot}$ and $\left[\widehat{Q}_i \right]_{i=1}^n = \left[\overline{Q}_i \right]_{i=1}^n$. According to Eq. (49), the constraints of advantage-based IGM stated in Fact 1 (Eq. (4), Eq. (5), and Eq. (7)) are satisfied, which means that $\left(\widehat{Q}_{tot}, \left[\widehat{Q}_i \right]_{i=1}^n \right) \in \widehat{\mathcal{Q}}$ and $\overline{\mathcal{Q}} \subseteq \widehat{\mathcal{Q}}$.

Thus, when assuming neural networks provide universal function approximation, the joint action-value function class that QPLEX can realize is equivalent to what is induced by the IGM principle. \square

B IMPLEMENTATION DETAILS AND EXPERIMENT SETTINGS

B.1 IMPLEMENTATION DETAILS

We adopt the PyMARL (Samvelyan et al., 2019) implementation of state-of-the-art baselines: QTRAN (Son et al., 2019), QMIX (Rashid et al., 2018), VDN (Sunehag et al., 2018), and Qatten (Yang et al., 2020). The hyper-parameters of these algorithms are the same as that in SMAC (Samvelyan et al., 2019) and referred in their source codes. QPLEX is also based on PyMARL, whose special hyper-parameters are illustrated in Table 1 and other common hyper-parameters are adopted by the default implementation of PyMARL (Samvelyan et al., 2019). Especially, in the online data collection, we take the advanced implementation of *Transformation* of Qatten in QPLEX. To ensure the positivity of important weights of *Transformation* and joint advantage function, we add a sufficiently small amount $\epsilon' = 10^{-10}$ on $[w_i]_{i=1}^n$ and $[\lambda_i]_{i=1}^n$. In addition, we stop gradients of local advantage function A_i to increase the optimization stability of the max operator of dueling structure. This instability consideration about max operator has been justified by Dueling DQN (Wang et al., 2016). We approximate the joint action-value function as

$$Q_{tot}(\tau, \mathbf{a}) \approx \sum_{i=1}^n Q_i(\tau, a_i) + \sum_{i=1}^n (\lambda_i(\tau, \mathbf{a}) - 1) \tilde{A}_i(\tau, a_i), \quad (50)$$

where \tilde{A}_i denotes a variant of the local advantage function A_i by stopping gradients.

Our training time on an NVIDIA RTX 2080TI GPU of each task is about 6 hours to 20 hours, depending on the agent number and the episode length limit of each map. The percentage of episodes in which MARL agents defeat all enemy units within the time limit is called *test win rate*. We pause training every 10k timesteps and evaluate 32 episodes with decentralized greedy action selection to measure *test win rate* of each algorithm. After training every 200 episodes, the target network will be updated once. We call this update period an *Iteration* for didactic tasks. In the two-state MMDP, *Optimal* line of Figure 3 is approximately $\sum_{i=0}^{99} \gamma^i = 63.4$ in one episode of 100 timesteps.

Training with Online Data Collection We have collected a total of 2 million timestep data for each task and test the model every 10 thousand steps. We use ϵ -greedy exploration and a limited

Map Name	Replay Buffer Size	Behaviour Test Win Rate	Behaviour Policy
2s3z	20k episodes	95.8%	QMIX
3s5z	20k episodes	92.0%	QMIX
1c3s5z	20k episodes	90.2%	QMIX
2s_vs_1sc	20k episodes	98.1%	QMIX
3s_vs_5z	20k episodes	94.4%	VDN
2c_vs_64zg	50k episodes	80.9%	QMIX

Table 2: The dataset configurations of offline data collection setting.

first-in-first-out (FIFO) replay buffer of size 5000 episodes, where ϵ is linearly annealed from 1.0 to 0.05 over 50k timesteps and keep it constant for the rest training process. To utilize the training buffer more efficiently, we perform gradient updates twice with a batch of 32 episodes after collecting every episode for each algorithm.

Training with Offline Data Collection To construct a diverse dataset, we train a behavior policy of QMIX (Rashid et al., 2018) or VDN (Sunehag et al., 2018) and collect its 20k or 50k experienced episodes throughout the training process. The dataset configurations are shown in Table 2. We evaluate QPLEX and four baselines over six random seeds, which includes three different datasets and tests two seeds on each dataset. We train 300 epochs to demonstrate our learning performance, where each epoch trains 160k transitions with a batch of 32 episodes. Moreover, the training process of behavior policy is the same as that discussed in PyMARL (Samvelyan et al., 2019).

B.2 STARCRAFT II

We consider the combat scenario of StarCraft II unit micromanagement tasks, where the enemy units are controlled by the built-in AI, and each ally unit is controlled by the reinforcement learning agent. The units of the two groups can be asymmetric but the units of each group should belong to the same race. At each timestep, every agent takes an action from the discrete action space, which includes the following actions: noop, move [direction], attack [enemy id], and stop. Under the control of these actions, agents move and attack in continuous maps. At each time step, MARL agents will get a global reward equal to the total damage done to enemy units. Killing each enemy unit and winning the combat will bring additional bonuses of 10 and 200, respectively. We briefly introduce the SMAC challenges of our paper in Table 3.

Map Name	Ally Units	Enemy Units
2s3z	2 Stalkers & 3 Zealots	2 Stalkers & 3 Zealots
3s5z	3 Stalkers & 5 Zealots	3 Stalkers & 5 Zealots
1c3s5z	1 Colossus, 3 Stalkers & 5 Zealots	1 Colossus, 3 Stalkers & 5 Zealots
1c3s8z_vs_1c3s9z	1 Colossus, 3 Stalkers & 8 Zealots	1 Colossus, 3 Stalkers & 9 Zealots
7sz	7 Stalkers & 7 Zealots	7 Stalkers & 7 Zealots
5s10z	5 Stalkers & 10 Zealots	5 Stalkers & 10 Zealots
2s_vs_1sc	2 Stalkers	1 Spine Crawler
3s_vs_5z	3 Stalkers	5 Zealots
2c_vs_64zg	2 Colossi	64 Zerglings

Table 3: SMAC challenges.

C DEFERRED FIGURES AND TABLES IN SECTION 4

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	8	-12	-12
$\mathcal{A}^{(2)}$	-12	0	0
$\mathcal{A}^{(3)}$	-12	0	0

(a) Payoff of matrix game

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	8.0	-12.1	-12.1
$\mathcal{A}^{(2)}$	-12.2	-0.0	-0.0
$\mathcal{A}^{(3)}$	-12.1	-0.0	-0.0

(b) Q_{tot} of QPLEX

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	8.0	-12.0	-12.0
$\mathcal{A}^{(2)}$	-12.0	-0.0	0.0
$\mathcal{A}^{(3)}$	-12.0	0.0	0.0

(c) Q_{tot} of QTRAN

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	-7.8	-7.8	-7.8
$\mathcal{A}^{(2)}$	-7.8	-0.0	-0.0
$\mathcal{A}^{(3)}$	-7.8	-0.0	-0.0

(d) Q_{tot} of QMIX

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	-6.4	-5.0	-5.0
$\mathcal{A}^{(2)}$	-5.0	-3.5	-3.5
$\mathcal{A}^{(3)}$	-5.0	-3.5	-3.5

(e) Q_{tot} of VDN

$a_2 \backslash a_1$	$\mathcal{A}^{(1)}$	$\mathcal{A}^{(2)}$	$\mathcal{A}^{(3)}$
$\mathcal{A}^{(1)}$	-6.5	-4.9	-4.9
$\mathcal{A}^{(2)}$	-5.0	-3.5	-3.4
$\mathcal{A}^{(3)}$	-5.0	-3.5	-3.5

(f) Q_{tot} of Qatten

Table 4: (a) Payoff matrix of the one-step game. Boldface means the optimal joint action selection from payoff matrix. (b-f) Deferred joint action-value functions Q_{tot} of QPLEX, QTRAN, QMIX, VDN, and Qatten. Boldface means greedy joint action selection from joint action-value functions.

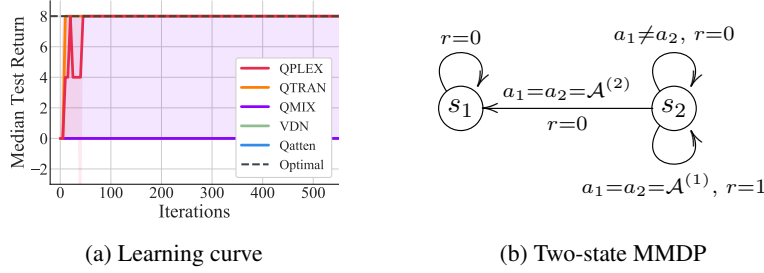


Figure 6: (a) The learning curves of QPLEX and other baselines. (b) A special two-state MMDP used to demonstrate the training stability of the multi-agent Q-learning algorithms. r is a shorthand for $r(s, \mathbf{a})$.

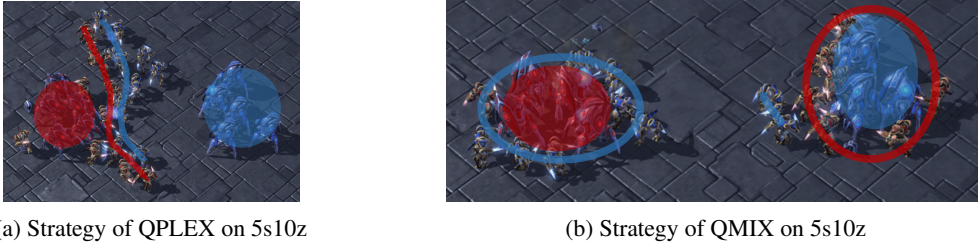


Figure 7: Visualized strategies of QPLEX and QMIX on 5s10z map of StarCraft II benchmark. Red marks represent learning agents, and blue marks represent build-in AI agents.

As shown in Figure 7, both MARL agents and opponents contain 5 ranged soldiers (denoted by a circle) and 10 melee soldiers (denoted by line) on 5s10z map. The ranged soldiers have stronger combat capabilities and need to be protected strategically. QPLEX uses 10 melee soldiers to build lines of defense against the enemy, while QMIX fails to coordinate melee soldiers such that ranged soldiers have to fight against the enemy directly.